

GLAD package: Gain and Loss Analysis of DNA

Philippe Hupe^{1,2} and Emmanuel Barillot²

May 9, 2009

1. UMR 144 CNRS/Institut Curie, Institut Curie, 26, rue d'Ulm, Paris, 75248 cedex 05, France
2. Service Bioinformatique, Institut Curie, 26, rue d'Ulm, Paris, 75248 cedex 05, France
`glad@curie.fr` <http://bioinfo.curie.fr>

Contents

1	Overview	2
2	Data	2
2.1	Public data set	2
2.2	Bladder cancer data	2
3	GLAD classes	2
3.1	arrayCGH	2
3.2	profileCGH and profileChr	2
4	Analysis of array CGH profile	3
4.1	Segmentation algorithms	3
4.2	The <i>glad</i> function	3
4.3	The <i>daglad</i> function	8
4.4	Tuning parameters	11
5	Graphical functions	11
5.1	Plot of raw array data	11
5.2	Plot of genomic profile	14

1 Overview

This document presents an overview of the GLAD package (Gain and Loss Analysis of DNA). This package is devoted to the analysis of Array Comparative Genomic Hybridization (array CGH) (Pinkel et al., 1998; Snijders et al., 2001; Solinas-Toldo et al., 1997; Ishkanian et al., 2004). The methodology for detecting the breakpoints delimiting altered regions in genomic patterns and assigning a status (normal, gained or lost) to each chromosomal region described in the paper Hupé et al. (2004) is implemented in this package. Some graphical functions are provided as well.

2 Data

2.1 Public data set

We used the public data set described in Snijders et al. (2001). The data corresponds to 15 human cell strains with known karyotypes (12 fibroblast cell strains, 2 chorionic villus cell strains, 1 lymphoblast cell strain) from the NIGMS Human Genetics Cell Repository (<http://locus.umdj.edu/nigms>). Each cell strain has been hybridized on CGH arrays of 2276 BACs, spotted in triplicates. Two array CGH profiles from the data obtained by Veltman et al. (2003) are available.

2.2 Bladder cancer data

Bladder cancer data from tumors collected at Henri Mondor Hospital (CrÃ©teil, France) (Billerey et al., 2001) have been hybridized on CGH arrays composed of 2464 BACs (Radvanyi, Pinkel et al., unpublished results). In this data, only the log-ratios are provided and no information about clones is available since the data is not yet published. This data allows only some graphical functionalities to be shown and will be used as a support to illustrate some functions for array normalization (not yet available in the current version of the package).

3 GLAD classes

3.1 arrayCGH

This class stores raw values after images analysis. The object arrayCGH is a list with at least a data.frame named arrayValues and a vector named arrayDesign. The data.frame arrayValues must contain the following fields:

Col Vector of columns coordinates.

Row Vector of rows coordinates.

... Other elements can be added.

The vector arrayDesign is composed of 4 values: c(arrayCol, arrayRow, SpotCol, SpotRow). The array CGH is represented by arrayRow*arrayCol blocs and each bloc is composed of SpotRow*SpotCol spots. N.B.: Col takes the values in 1:arrayRow*SpotRow and Row takes the values in 1:arrayCol*SpotCol

3.2 profileCGH and profileChr

This class stores synthetic values related to each clone available on the arrayCGH. The object profileChr corresponds to data of only one chromosome. Objects profileCGH and profileChr are composed of a list with the first element profileValues which is a data.frame with the following columns names:

LogRatio Test over Reference log-ratio.

PosOrder The rank position of each clone on the genome.

PosBase The base position of each clone on the genome.

Chromosome Chromosome name.

Clone The name of the corresponding clone.

... Other elements can be added.

LogRatio, Chromosome and PosOrder are compulsory.

To create those objects you can use the function *as.profileCGH*.

4 Analysis of array CGH profile

Two functions are available: *glad* and *daglad*. The second one is an improvement of the first one which was originally described in Hupé et al. (2004). We recommend to use the *daglad* function. For fast computation use the option *smoothfunc=haarseg*.

4.1 Segmentation algorithms

Two algorithms are available for data segmentation:

- AWS (Polzehl and Spokoiny, 2000, 2002)
- HaarSeg (Ben-Yaacov and Eldar, 2008)

4.2 The *glad* function

A result of the GLAD methodology on cell line gm13330 (Snijders et al., 2001) is presented in **Figure 1**.

```
> data(snijders)
> profileCGH <- as.profileCGH(gm13330)
> res <- glad(profileCGH, mediancenter = FALSE, smoothfunc = "lawsglad",
+   bandwidth = 10, round = 1.5, model = "Gaussian", lkern = "Exponential",
+   qlambda = 0.999, base = FALSE, lambdabreak = 8, lambdacluster = 8,
+   lambdaclusterGen = 40, type = "tricubic", param = c(d = 6),
+   alpha = 0.001, msize = 5, method = "centroid", nmax = 8,
+   verbose = FALSE)

[1] "Smoothing for each Chromosome"
[1] "Optimization of the Breakpoints"
[1] "subsetdata"
      user  system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.001899958 secs
[1] "Temps findCluster: 0.0846970081329346"
[1] "subsetdata"
      user  system elapsed
```

```

0.004 0.000 0.001
[1] "diff"
Time difference of 0.001661062 secs
[1] "Temps findCluster: 0.0210549831390381"
[1] "subsetdata"
    user system elapsed
0.000 0.000 0.001
[1] "diff"
Time difference of 0.001652002 secs
[1] "Temps findCluster: 0.0214431285858154"
[1] "subsetdata"
    user system elapsed
    0      0      0
[1] "diff"
Time difference of 0.001733065 secs
[1] "Temps findCluster: 0.0838251113891602"
[1] "subsetdata"
    user system elapsed
0.000 0.000 0.001
[1] "diff"
Time difference of 0.001657009 secs
[1] "Temps findCluster: 0.0209591388702393"
[1] "subsetdata"
    user system elapsed
0.000 0.000 0.001
[1] "diff"
Time difference of 0.001670122 secs
[1] "Temps findCluster: 0.0215580463409424"
[1] "subsetdata"
    user system elapsed
    0      0      0
[1] "diff"
Time difference of 0.001676083 secs
[1] "Temps findCluster: 0.0210118293762207"
[1] "subsetdata"
    user system elapsed
0.000 0.000 0.001
[1] "diff"
Time difference of 0.001661062 secs
[1] "Temps findCluster: 0.0216858386993408"
[1] "subsetdata"
    user system elapsed
    0      0      0
[1] "diff"
Time difference of 0.001654863 secs
[1] "Temps findCluster: 0.0210039615631104"
[1] "subsetdata"
    user system elapsed
    0      0      0
[1] "diff"

```

```

Time difference of 0.001658916 secs
[1] "Temps findCluster: 0.0216009616851807"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001678944 secs
[1] "Temps findCluster: 0.0210578441619873"
[1] "subsetdata"
      user  system elapsed
0.000 0.000 0.001
[1] "diff"
Time difference of 0.001662016 secs
[1] "Temps findCluster: 0.0216538906097412"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001655102 secs
[1] "Temps findCluster: 0.021273136138916"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001657009 secs
[1] "Temps findCluster: 0.0215368270874023"
[1] "subsetdata"
      user  system elapsed
0.000 0.000 0.001
[1] "diff"
Time difference of 0.001663923 secs
[1] "Temps findCluster: 0.0211300849914551"
[1] "subsetdata"
      user  system elapsed
0.000 0.000 0.001
[1] "diff"
Time difference of 0.00177598 secs
[1] "Temps findCluster: 0.0216610431671143"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001650095 secs
[1] "Temps findCluster: 0.0210628509521484"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001657963 secs
[1] "Temps findCluster: 0.0215399265289307"

```

```

[1] "subsetdata"
      user system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.001649141 secs
[1] "Temps findCluster: 0.0209999084472656"
[1] "subsetdata"
      user system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.001778126 secs
[1] "Temps findCluster: 0.0215849876403809"
[1] "subsetdata"
      user system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001665115 secs
[1] "Temps findCluster: 0.0210440158843994"
[1] "subsetdata"
      user system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001673937 secs
[1] "Temps findCluster: 0.0214419364929199"
[1] "subsetdata"
      user system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001660109 secs
[1] "Temps findCluster: 0.0211269855499268"
[1] "subsetdata"
      user system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.003628016 secs
[1] "Temps findCluster: 0.162232160568237"
[1] "Results Preparation"

```

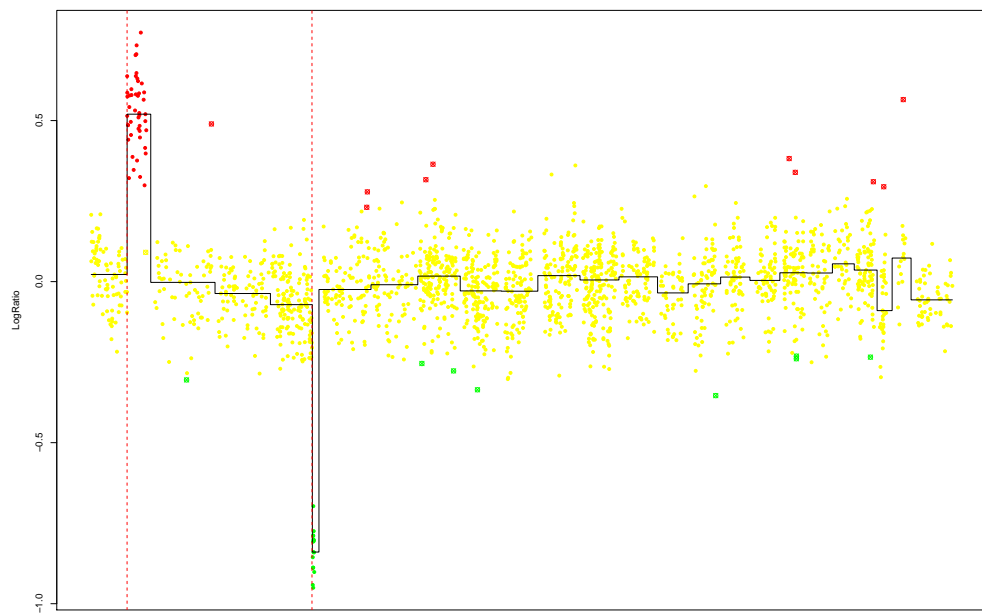


Figure 1: Results of glad on cell line gm13330 (Snijders data).

4.3 The *daglad* function

The algorithm implemented in this function is a slightly modified version of the GLAD algorithm.

```
> data(veltman)
> profileCGH <- as.profileCGH(P9)
> profileCGH <- daglad(profileCGH, mediancenter = FALSE, normalrefcenter = FALSE,
+   genomestep = FALSE, smoothfunc = "lawsglad", lkern = "Exponential",
+   model = "Gaussian", qlambda = 0.999, bandwidth = 10, base = FALSE,
+   round = 1.5, lambdabreak = 8, lambdaclusterGen = 40, param = c(d = 6),
+   alpha = 0.001, msize = 5, method = "centroid", nmin = 1,
+   nmax = 8, amplicon = 1, deletion = -5, deltaN = 0.2, forceGL = c(-0.3,
+   0.3), nbsigma = 3, MinBkpWeight = 0.35, CheckBkpPos = TRUE)

[1] "Smoothing for each Chromosome"
[1] "Optimization of the Breakpoints"
[1] "DNA copy number calling"
[1] "subsetdata"
      user  system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.002655983 secs
[1] "Temps findCluster: 0.165307998657227"
[1] "jointure BkpInfo"
      user  system elapsed
0.000   0.000   0.002
[1] "Check Breakpoints Position"
[1] "Results Preparation"
```

The *daglad* function allows to choose some threshold to help the algorithm to identify the status of the genomic regions. The thresholds are given in the following parameters:

- deltaN
- forceGL
- deletion
- amplicon

Comparing **Figure 2** and **Figure 3** shows the influence of two different sets of parameters.

```
> data(veltman)
> profileCGH <- as.profileCGH(P9)
> profileCGH <- daglad(profileCGH, mediancenter = FALSE, normalrefcenter = FALSE,
+   genomestep = FALSE, smoothfunc = "lawsglad", lkern = "Exponential",
+   model = "Gaussian", qlambda = 0.999, bandwidth = 10, base = FALSE,
+   round = 1.5, lambdabreak = 8, lambdaclusterGen = 40, param = c(d = 6),
+   alpha = 0.001, msize = 5, method = "centroid", nmin = 1,
+   nmax = 8, amplicon = 1, deletion = -5, deltaN = 0.1, forceGL = c(-0.15,
+   0.15), nbsigma = 3, MinBkpWeight = 0.35, CheckBkpPos = TRUE)
```

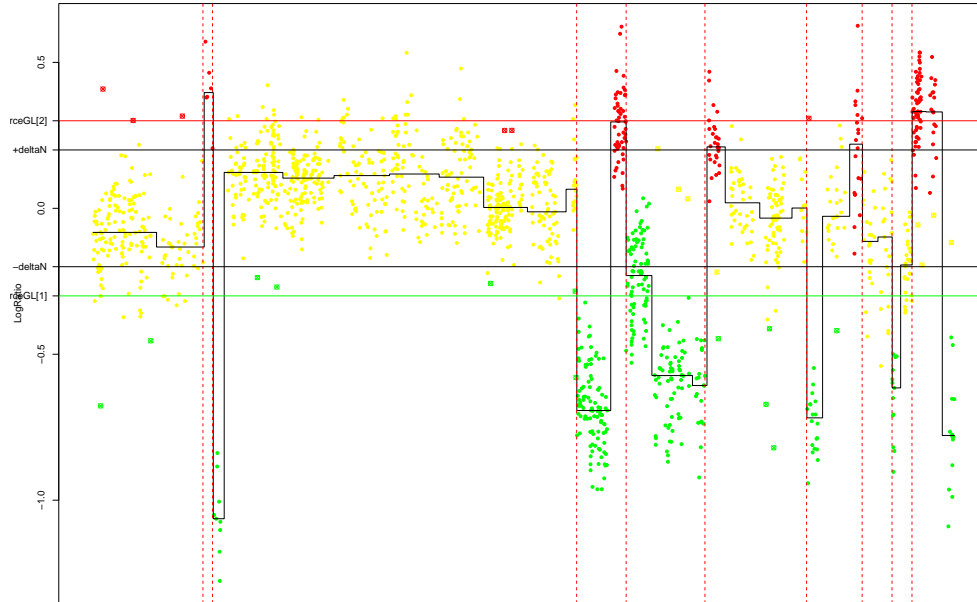



Figure 2: Results of daglad on the patient P9 (Veltman data).

```
[1] "Smoothing for each Chromosome"
[1] "Optimization of the Breakpoints"
[1] "DNA copy number calling"
[1] "subsetdata"
      user system elapsed
0.004   0.000   0.000
[1] "diff"
Time difference of 0.003164053 secs
[1] "Temps findCluster: 0.166804790496826"
[1] "jointure BkpInfo"
      user system elapsed
0.000   0.000   0.001
[1] "Check Breakpoints Position"
[1] "Results Preparation"
```

The *daglad* function allows a smoothing step over the whole genome (if *genomestep=TRUE*) where all the chromosomes are concatenated together. During this step, the cluster which corresponds to the Normal DNA level is identified: the thresholds used in the function (*deltaN*, *forceGL*, *amplicon*, *deletion*) are then compared to the median of this cluster.

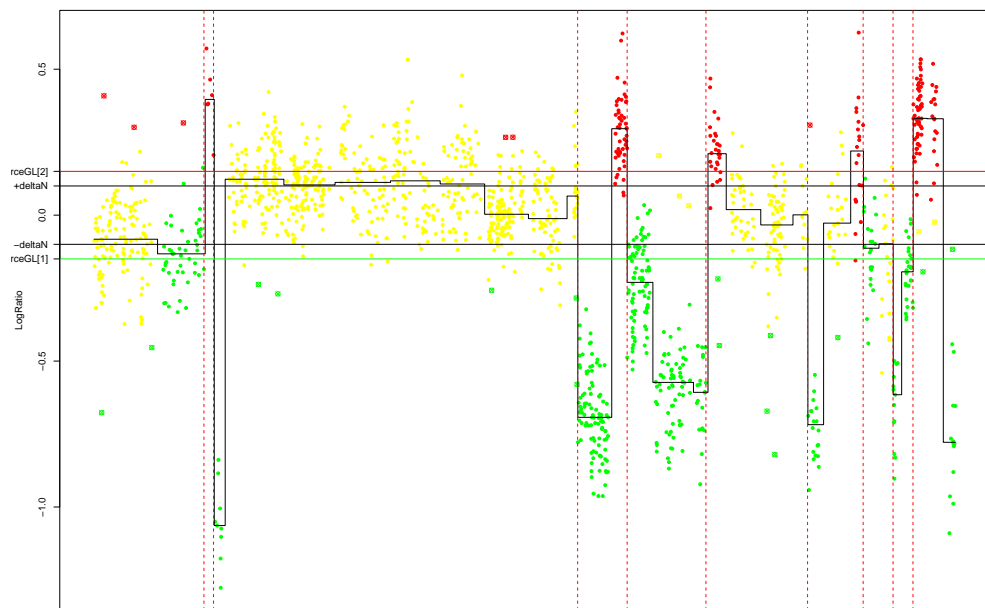


Figure 3: Results of daglad on the patient P9 (Veltman data) - Influence of the thresholds.

4.4 Tuning parameters

The most important parameters are:

- *lambdabreak*
- *lambdacluster*
- *lambdaclusterGen*
- *param* $c(d = 6)$

Decreasing those parameters will lead to a higher number of breakpoints identified. For arrays experiments with very small Signal to Noise ratio it is recommended to use a small value of *param* like $d = 3$ or less.

5 Graphical functions

5.1 Plot of raw array data

```
> data(arrayCGH)
> array <- list(arrayValues = array2, arrayDesign = c(4, 4, 21,
+ 22))
> class(array) <- "arrayCGH"
```

```
> arrayPlot(array, "Log2Rat", bar = "none")
```

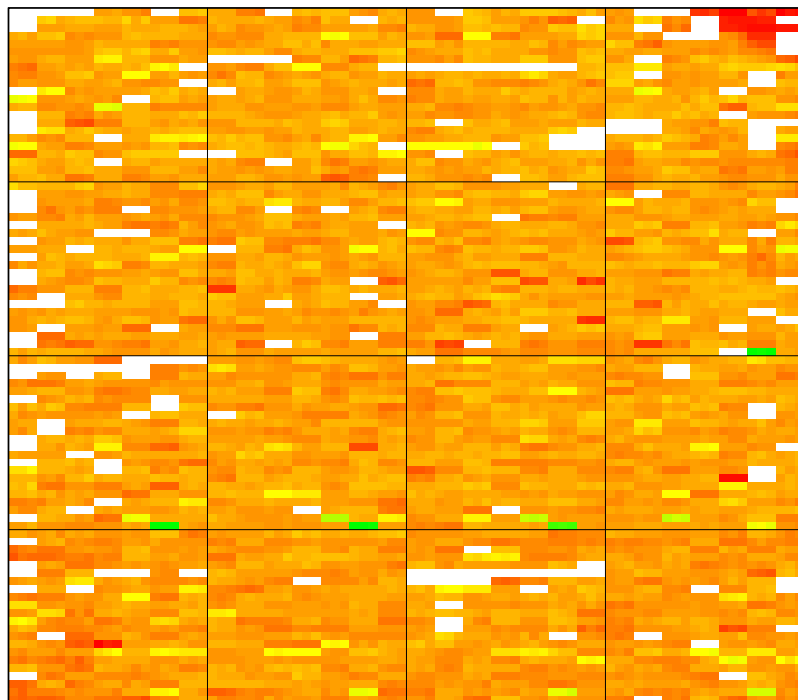


Figure 4: Spatial image of array CGH

```
> arrayPersp(array, "Log2Rat", box = FALSE, theta = 110, phi = 40,  
+           bar = FALSE)
```

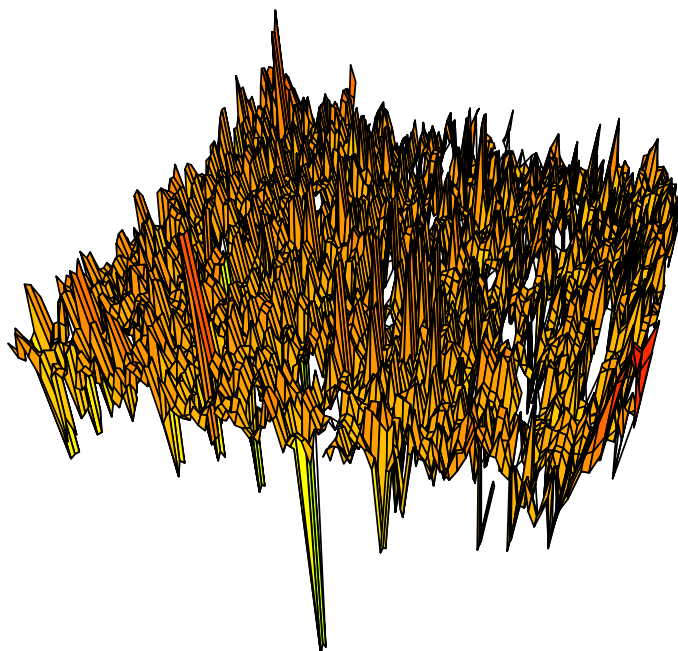


Figure 5: Perspective image of array CGH

5.2 Plot of genomic profile

```
[1] "Smoothing for each Chromosome"
[1] "Optimization of the Breakpoints"
[1] "subsetdata"
      user system elapsed
0.000  0.000  0.001
[1] "diff"
Time difference of 0.001725912 secs
[1] "Temps findCluster: 0.0949280261993408"
[1] "subsetdata"
      user system elapsed
0.000  0.000  0.001
[1] "diff"
Time difference of 0.001639128 secs
[1] "Temps findCluster: 0.0242700576782227"
[1] "subsetdata"
      user system elapsed
0.004  0.000  0.000
[1] "diff"
Time difference of 0.001636028 secs
[1] "Temps findCluster: 0.0238020420074463"
[1] "subsetdata"
      user system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001713991 secs
[1] "Temps findCluster: 0.0942051410675049"
[1] "subsetdata"
      user system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001631975 secs
[1] "Temps findCluster: 0.0236990451812744"
[1] "subsetdata"
      user system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001632929 secs
[1] "Temps findCluster: 0.0238149166107178"
[1] "subsetdata"
      user system elapsed
0.004  0.000  0.000
[1] "diff"
Time difference of 0.001638889 secs
[1] "Temps findCluster: 0.0240628719329834"
[1] "subsetdata"
      user system elapsed
      0      0      0
[1] "diff"
```

```

Time difference of 0.001636982 secs
[1] "Temps findCluster: 0.0238702297210693"
[1] "subsetdata"
      user system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.001744986 secs
[1] "Temps findCluster: 0.0247199535369873"
[1] "subsetdata"
      user system elapsed
0         0         0
[1] "diff"
Time difference of 0.001851082 secs
[1] "Temps findCluster: 0.0240039825439453"
[1] "subsetdata"
      user system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.001912117 secs
[1] "Temps findCluster: 0.117748022079468"
[1] "subsetdata"
      user system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.001662970 secs
[1] "Temps findCluster: 0.0260658264160156"
[1] "subsetdata"
      user system elapsed
0.000   0.000   0.001
[1] "diff"
Time difference of 0.001660109 secs
[1] "Temps findCluster: 0.0242810249328613"
[1] "subsetdata"
      user system elapsed
0         0         0
[1] "diff"
Time difference of 0.001863003 secs
[1] "Temps findCluster: 0.0243391990661621"
[1] "subsetdata"
      user system elapsed
0         0         0
[1] "diff"
Time difference of 0.00166893 secs
[1] "Temps findCluster: 0.0242240428924561"
[1] "subsetdata"
      user system elapsed
0         0         0
[1] "diff"
Time difference of 0.001670122 secs
[1] "Temps findCluster: 0.0244688987731934"

```

```

[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001757145 secs
[1] "Temps findCluster: 0.0240330696105957"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001643896 secs
[1] "Temps findCluster: 0.0241880416870117"
[1] "subsetdata"
      user  system elapsed
    0.000  0.000  0.001
[1] "diff"
Time difference of 0.001653910 secs
[1] "Temps findCluster: 0.0243849754333496"
[1] "subsetdata"
      user  system elapsed
    0.004  0.000  0.000
[1] "diff"
Time difference of 0.001752853 secs
[1] "Temps findCluster: 0.0240609645843506"
[1] "subsetdata"
      user  system elapsed
    0.000  0.000  0.001
[1] "diff"
Time difference of 0.001640081 secs
[1] "Temps findCluster: 0.0239520072937012"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.001642942 secs
[1] "Temps findCluster: 0.0242431163787842"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.00176692 secs
[1] "Temps findCluster: 0.0239670276641846"
[1] "subsetdata"
      user  system elapsed
      0      0      0
[1] "diff"
Time difference of 0.003554106 secs
[1] "Temps findCluster: 0.179353952407837"
[1] "Results Preparation"

```



```
> plotProfile(res, unit = 3, Bkp = TRUE, labels = FALSE, Smoothing = "Smoothing",
+   plotband = FALSE)
```

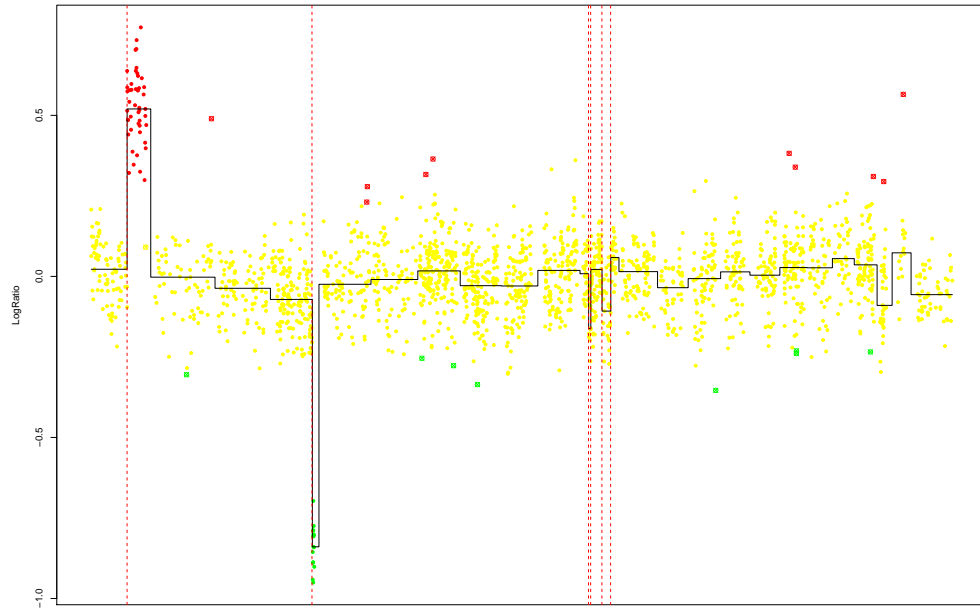


Figure 6: Genomic profile on the whole genome

```
> plotProfile(res, unit = 3, Bkp = TRUE, labels = FALSE, Smoothing = "Smoothing")
```

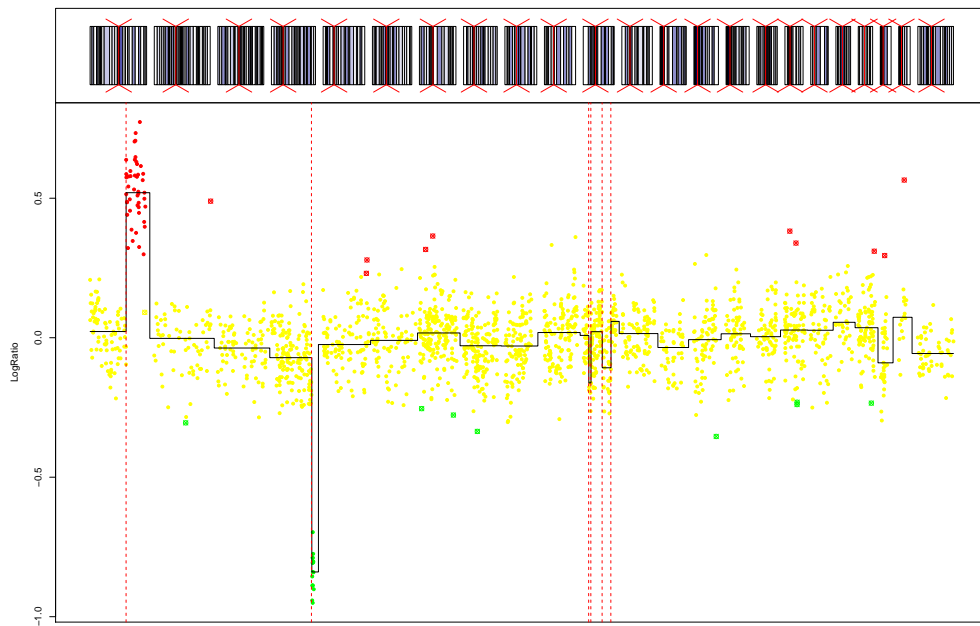


Figure 7: Genomic profile on the whole genome and cytogenetic banding

```

> text <- list(x = c(90000, 2e+05), y = c(0.15, 0.3), labels = c("NORMAL",
+   "GAIN"), cex = 2)
> plotProfile(res, unit = 3, Bkp = TRUE, labels = TRUE, Chromosome = 1,
+   Smoothing = "Smoothing", plotband = FALSE, text = text)

```

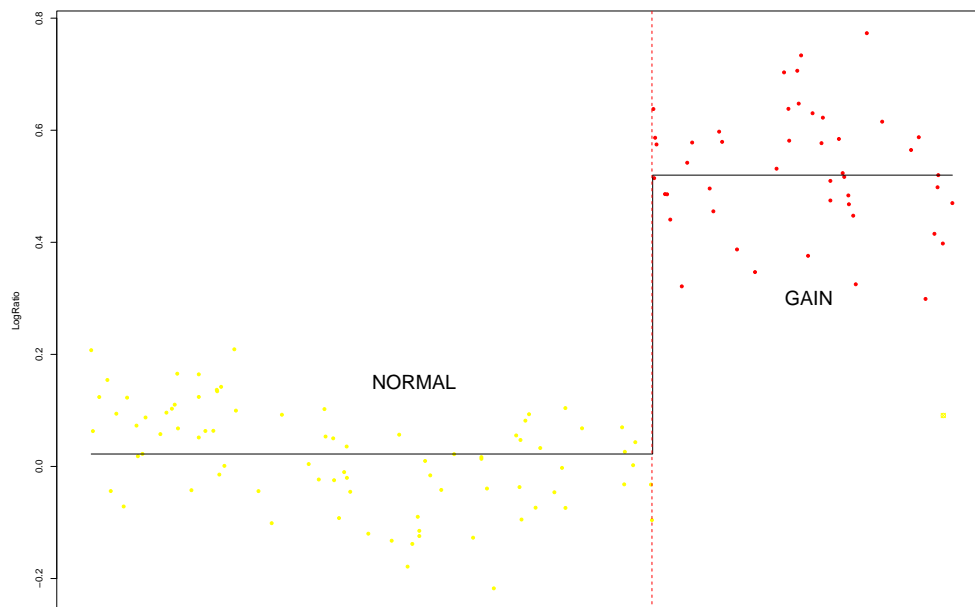


Figure 8: Genomic profile for chromosome 1

```

> text <- list(x = c(90000, 2e+05), y = c(0.15, 0.3), labels = c("NORMAL",
+   "GAIN"), cex = 2)
> plotProfile(res, unit = 3, Bkp = TRUE, labels = TRUE, Chromosome = 1,
+   Smoothing = "Smoothing", text = text, main = "Chromosome 1")

```

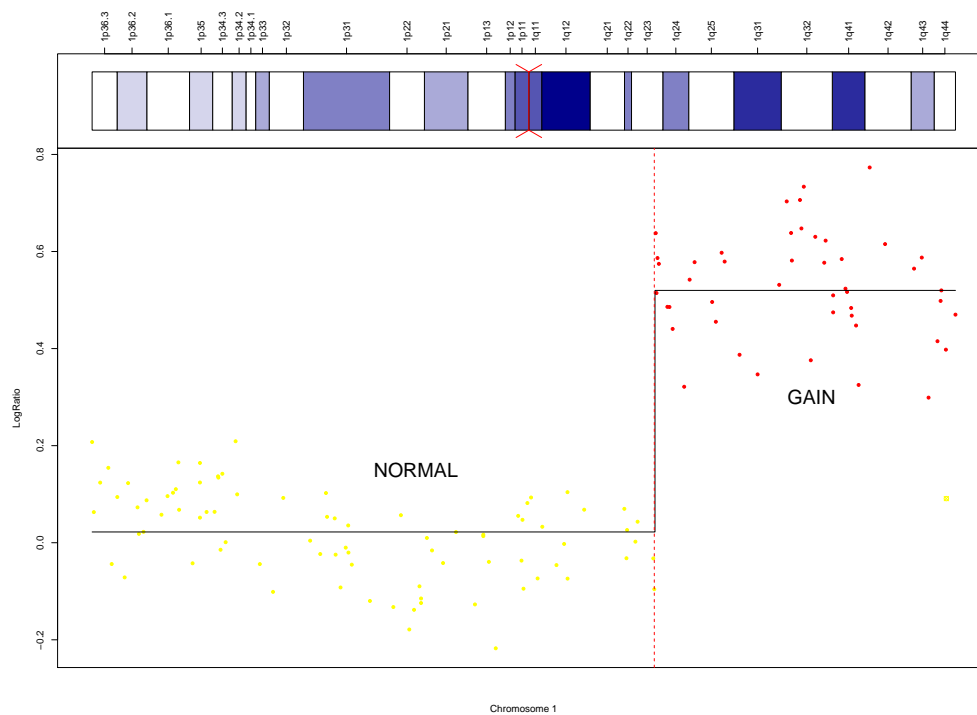


Figure 9: Genomic profile for chromosome 1 and cytogenetic banding with labels

References

- Ben-Yaacov, E. and Eldar, Y. C. (2008). A fast and flexible method for the segmentation of acgh data. *Bioinformatics*, 24:i139–i145.
- Billerey, C., Chopin, D., Aubriot-Lorton, M. H., Ricol, D., de Medina, S. G. D., Rhijn, B. V., Bralet, M. P., Lefrere-Belda, M. A., Lahaye, J. B., Abbou, C. C., Bonaventure, J., Zafrani, E. S., van der Kwast, T., Thiery, J. P., and Radvanyi, F. (2001). Frequent FGFR3 mutations in papillary non-invasive bladder (pTa) tumors. *Am. J. Pathol.*, 158:955–1959.
- Hupé, P., Stransky, N., Thiery, J. P., Radvanyi, F., and Barillot, E. (2004). Analysis of array CGH data: from signal ratio to gain and loss of DNA regions. *Bioinformatics*, 20:3413–3422.
- Ishkanian, A. S., Malloff, C. A., Watson, S. K., DeLeeuw, R. J., Chi, B., Coe, B. P., Snijders, A., Albertson, D. G., Pinkel, D., Marra, M. A., Ling, V., MacAulay, C., and Lam, W. L. (2004). A tiling resolution DNA microarray with complete coverage of the human genome. *Nat. Genet.*, 36:299–303.
- Pinkel, D., Segreaves, R., Sudar, D., Clark, S., Poole, I., Kowbel, D., Collins, C., Kuo, W. L., Chen, C., Zhai, Y., Dairkee, S. H., Ljung, B. M., Gray, J. W., and Albertson, D. G. (1998). High resolution analysis of dna copy number variation using comparative genomic hybridization to microarrays. *Nat. Genet.*, 20:207–211.
- Polzehl, J. and Spokoiny, S. (2000). Adaptive weights smoothing with applications to image restoration. *J. Roy. Statistical Society, Series B*, pages 335–354.
- Polzehl, J. and Spokoiny, S. (2002). Local likelihood modelling by adaptive weights smoothing. WIAS-Preprint 787.
- Snijders, A. M., Nowak, N., Segreaves, R., Blackwood, S., Brown, N., Conroy, J., Hamilton, G., Hindle, A. K., Huey, B., Kimura, K., S, S. L., Myambo, K., Palmer, J., Ylstra, B., Yue, J. P., Gray, J. W., Jain, A. N., Pinkel, D., and Albertson, D. G. (2001). Assembly of microarrays for genome-wide measurement of dna copy number. *Nat. Genet.*, 29:263–4.
- Solinas-Toldo, S., Lampel, S., Stilgenbauer, S., Nickolenko, J., Benner, A., Dohner, H., Cremer, T., and Lichter, P. (1997). Matrix-based comparative genomic hybridization: Biochips to screen for genomic imbalances. *Genes Chromosomes Cancer*, 20:399–407.
- Veltman, J. A., Fridlyand, J., Pejavar, S., Olshen, A. B., Korkola, J. E., DeVries, S., Carroll, P., Kuo, W.-L., Pinkel, D., Albertson, D., Cordon-Cardo, C., Jain, A. N., and Waldman, F. M. (2003). Array-based comparative genomic hybridization for genome-wide screening of DNA copy number in bladder tumors. *Cancer Res*, 63:2872–2880.